

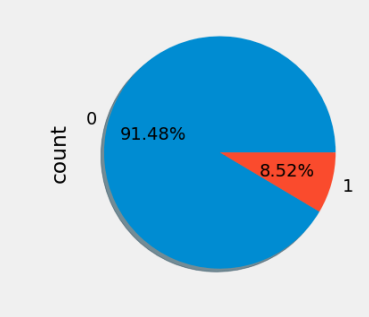
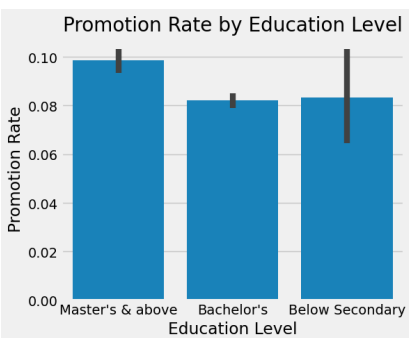
Data Collection and Preprocessing Phase

Date	12 June 2025
Team ID	SWTID1749627644
Project Title	Human Resource Management: Predicting Employee Promotions using Machine Learning
Maximum Marks	6 Marks

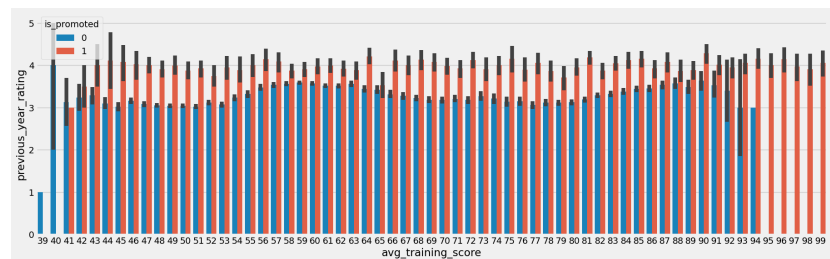
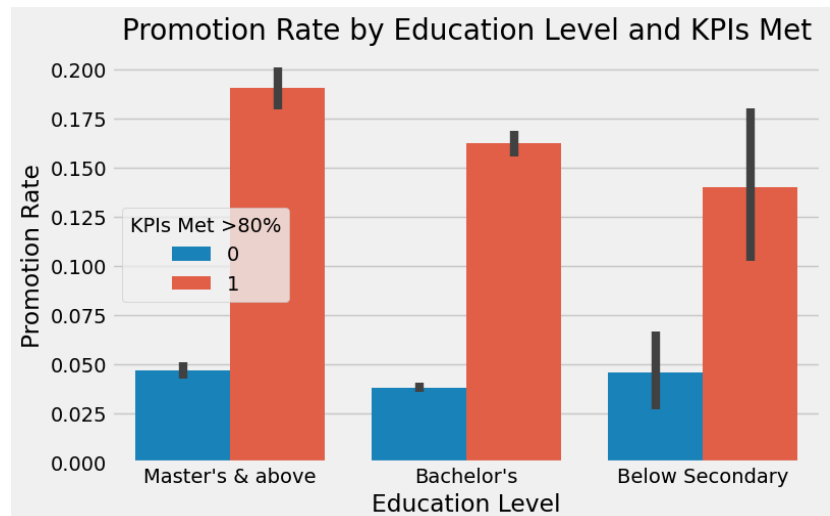
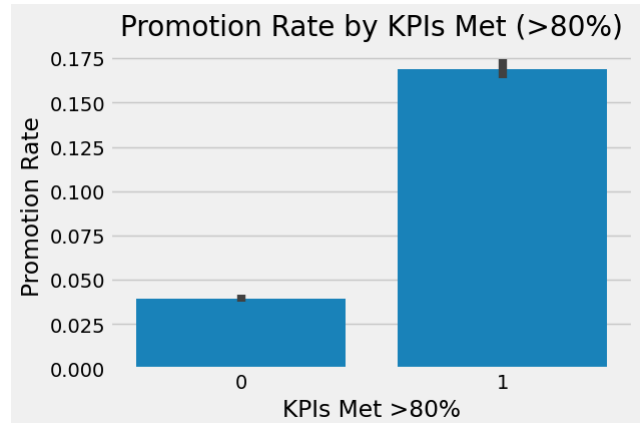
Data Exploration and Preprocessing:

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

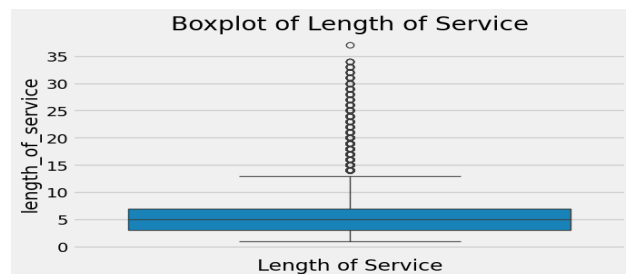
Section	Description																																																						
Data Overview	<u>Dimension:</u> 54808 rows x 14 columns																																																						
	<u>Descriptive Statistics:</u>																																																						
	<table><tr><th></th><th>employee_id</th><th>no_of_trainings</th><th>age</th><th>previous_year_rating</th><th>length_of_service</th></tr><tr><td>count</td><td>54808.000000</td><td>54808.000000</td><td>54808.000000</td><td>50684.000000</td><td>54808.000000</td></tr><tr><td>mean</td><td>39195.830627</td><td>1.253011</td><td>34.803915</td><td>3.329256</td><td>5.865512</td></tr><tr><td>std</td><td>22586.581449</td><td>0.609264</td><td>7.660169</td><td>1.259993</td><td>4.265094</td></tr><tr><td>min</td><td>1.000000</td><td>1.000000</td><td>20.000000</td><td>1.000000</td><td>1.000000</td></tr><tr><td>25%</td><td>19669.750000</td><td>1.000000</td><td>29.000000</td><td>3.000000</td><td>3.000000</td></tr><tr><td>50%</td><td>39225.500000</td><td>1.000000</td><td>33.000000</td><td>3.000000</td><td>5.000000</td></tr><tr><td>75%</td><td>58730.500000</td><td>1.000000</td><td>39.000000</td><td>4.000000</td><td>7.000000</td></tr><tr><td>max</td><td>78298.000000</td><td>10.000000</td><td>60.000000</td><td>5.000000</td><td>37.000000</td></tr></table>		employee_id	no_of_trainings	age	previous_year_rating	length_of_service	count	54808.000000	54808.000000	54808.000000	50684.000000	54808.000000	mean	39195.830627	1.253011	34.803915	3.329256	5.865512	std	22586.581449	0.609264	7.660169	1.259993	4.265094	min	1.000000	1.000000	20.000000	1.000000	1.000000	25%	19669.750000	1.000000	29.000000	3.000000	3.000000	50%	39225.500000	1.000000	33.000000	3.000000	5.000000	75%	58730.500000	1.000000	39.000000	4.000000	7.000000	max	78298.000000	10.000000	60.000000	5.000000	37.000000
		employee_id	no_of_trainings	age	previous_year_rating	length_of_service																																																	
	count	54808.000000	54808.000000	54808.000000	50684.000000	54808.000000																																																	
	mean	39195.830627	1.253011	34.803915	3.329256	5.865512																																																	
	std	22586.581449	0.609264	7.660169	1.259993	4.265094																																																	
	min	1.000000	1.000000	20.000000	1.000000	1.000000																																																	
	25%	19669.750000	1.000000	29.000000	3.000000	3.000000																																																	
	50%	39225.500000	1.000000	33.000000	3.000000	5.000000																																																	
75%	58730.500000	1.000000	39.000000	4.000000	7.000000																																																		
max	78298.000000	10.000000	60.000000	5.000000	37.000000																																																		

	<table><tr><th>KPIs_met >80%</th><th>awards_won?</th><th>avg_training_score</th><th>is_promoted</th></tr><tr><td>54808.000000</td><td>54808.000000</td><td>54808.000000</td><td>54808.000000</td></tr><tr><td>0.351974</td><td>0.023172</td><td>63.386750</td><td>0.085170</td></tr><tr><td>0.477590</td><td>0.150450</td><td>13.371559</td><td>0.279137</td></tr><tr><td>0.000000</td><td>0.000000</td><td>39.000000</td><td>0.000000</td></tr><tr><td>0.000000</td><td>0.000000</td><td>51.000000</td><td>0.000000</td></tr><tr><td>0.000000</td><td>0.000000</td><td>60.000000</td><td>0.000000</td></tr><tr><td>1.000000</td><td>0.000000</td><td>76.000000</td><td>0.000000</td></tr><tr><td>1.000000</td><td>1.000000</td><td>99.000000</td><td>1.000000</td></tr></table>	KPIs_met >80%	awards_won?	avg_training_score	is_promoted	54808.000000	54808.000000	54808.000000	54808.000000	0.351974	0.023172	63.386750	0.085170	0.477590	0.150450	13.371559	0.279137	0.000000	0.000000	39.000000	0.000000	0.000000	0.000000	51.000000	0.000000	0.000000	0.000000	60.000000	0.000000	1.000000	0.000000	76.000000	0.000000	1.000000	1.000000	99.000000	1.000000
KPIs_met >80%	awards_won?	avg_training_score	is_promoted																																		
54808.000000	54808.000000	54808.000000	54808.000000																																		
0.351974	0.023172	63.386750	0.085170																																		
0.477590	0.150450	13.371559	0.279137																																		
0.000000	0.000000	39.000000	0.000000																																		
0.000000	0.000000	51.000000	0.000000																																		
0.000000	0.000000	60.000000	0.000000																																		
1.000000	0.000000	76.000000	0.000000																																		
1.000000	1.000000	99.000000	1.000000																																		
Univariate Analysis	<div><div><p>Distribution of Average Training Score</p></div><div><p>Count of Education Levels</p></div></div> <div></div>																																				
Bivariate Analysis	<div><p>Promotion Rate by Education Level</p></div> <div><p>Training Score vs. Promotion</p></div>																																				

Multivariate Analysis



Outliers and Anomalies



	Number of outliers: 3489		
	length_of_service	is_promoted	
	13	16	0
	42	26	0
	60	17	1
	74	14	0
	99	17	0

	54691	19	0
	54695	18	1
	54697	15	0
	54754	14	0
	54803	17	0
	3489 rows × 2 columns		

Data Preprocessing Code Screenshots

Loading Data

```
df = pd.read_csv('../Dataset/emp_promotion.csv')  
print('Shape of train data {}'.format(df.shape))
```

✓ 0.1s

Shape of train data (54808, 14)

```
df.head()
```

✓ 0.0s

	employee_id	department	region	education	gender	recruitment_channel
0	65438	Sales & Marketing	region_7	Master's & above	f	sourcing
1	65141	Operations	region_22	Bachelor's	m	other
2	7513	Sales & Marketing	region_19	Bachelor's	m	sourcing
3	2542	Sales & Marketing	region_23	Bachelor's	m	other
4	48945	Technology	region_26	Bachelor's	m	other

Handling Missing Data

```
df.isnull().sum()
```

```
department      0
education      2409
no_of_trainings  0
age             0
previous_year_rating  4124
length_of_service  0
KPIs_met >80%    0
awards_won?     0
avg_training_score  0
is_promoted     0
dtype: int64
```

```
print(df['education'].value_counts())
df['education'] = df['education'].fillna(df['education'].mode()[0])
```

```
education
Bachelor's      36669
Master's & above 14925
Below Secondary   805
Name: count, dtype: int64
```

```
print(df['previous_year_rating'].value_counts())
df['previous_year_rating'] = df['previous_year_rating'].fillna(df['previous_year_rating'].mode()[0])

previous_year_rating
3.0    18618
5.0    11741
4.0     9877
1.0     6223
2.0     4225
Name: count, dtype: int64
```

Data Transformation

We do not need employee id, gender, region, recruitment channel attributes for predicting promotion, so we will be removing these unwanted features.

```
df = df.drop(['employee_id', 'gender', 'region', 'recruitment_channel'], axis=1)
```

Python

Removing inconsistent rows

```
df.drop(index=[31860, 51374], inplace=True)
```

Capping outliers

```
df['length_of_service'] = [upperBound if x > upperBound else x for x in df['length_of_service']]
```

Feature Engineering

```
# Feature mapping on education column
df['education'] = df['education'].replace(("Below Secondary", "Bachelor's", "Master's & above"), (1,2,3))
```

	<pre>lb = LabelEncoder() df['department'] = lb.fit_transform(df['department']) sm = SMOTE() x_resample, y_resample = sm.fit_resample(x,y)</pre>
Save Processed Data	-